

Sample complexity of the distinct elements problem

Yihong Wu and Pengkun Yang*

December 13, 2016

Abstract

We consider the distinct elements problem, where the goal is to estimate the number of distinct colors in an urn containing k balls from repeated draws. We propose an estimator, based on sampling without replacement, with additive error guarantee. The sample complexity is optimal within $O(\log \log k)$ factors, and in fact within constant factors for most accuracy parameters. The optimal sample complexity is also applicable to sampling without replacement provided the sample size is a vanishing fraction of the urn size.

One of the key auxiliary results is a sharp bound on the minimum singular values of a real rectangular Vandermonde matrix, which might be of independent interest.

Keywords sampling large population, nonparametric statistics, discrete polynomial approximation, orthogonal polynomials, Vandermonde matrix, minimaxity

AMS 2010 subject classifications Primary: 62G05; secondary: 62C20, 62D05, 41A05, 41A10

Acknowledgement This research has been supported in part by the National Science Foundation under the grant agreement IIS-14-47879 and CCF-15-27105. The authors thank Greg Valiant for helpful discussion on Lemma 5.14 in his thesis [Val12].

*Y. Wu is with the Department of Statistics, Yale University, New Haven, CT, yihong.wu@yale.edu. P. Yang is with the Department of Electrical and Computer Engineering and the Coordinated Science Lab, University of Illinois at Urbana-Champaign, Urbana, IL, pyang14@illinois.edu.

Contents

| | | |
|----------|--|-----------|
| 1 | The Distinct Elements problem | 2 |
| 1.1 | Main results | 2 |
| 1.2 | Related work | 4 |
| 1.3 | Organization | 5 |
| 1.4 | Notations | 6 |
| 2 | Linear estimators via discrete polynomial approximation | 6 |
| 2.1 | Guarantees on general linear estimators | 6 |
| 2.2 | Exact solution to the ℓ_2 -approximation | 9 |
| 2.3 | Minimum singular values of real rectangle Vandermonde matrices | 12 |
| 2.4 | Lagrange interpolating polynomials and Stirling numbers | 15 |
| 3 | Optimality of the sample complexity | 17 |
| A | Connections between various sampling models | 21 |
| B | Correlation decay between fingerprints | 22 |
| C | Proof of auxiliary lemmas | 23 |
| D | Proof of results in Table 1 | 25 |

1 The Distinct Elements problem

The Distinct Elements problem [CCMN00] refers to the following question:

Given n balls randomly drawn from an urn containing k colored balls, how to estimate the total number of distinct colors in the urn?

Originating from ecology, numismatics and linguistics, this problem is also known as the *species problem* in the statistics literature [Lo92,BF93]. Apart from the theoretical interests, it has a wide array of applications in various fields, such as estimating the number of species in a population of animals [FCW43], the number of dies used to mint an ancient coinage [Est86], and the vocabulary size of an author [ET76]. In computer science, this problem frequently arises in large-scale databases, network monitoring, and data mining [RRSS09,BYJK⁺02,CCMN00], where the objective is to estimate the types of database entries or IP addresses from limited observations, since it is typically impossible to have full access to the entire database or keep track of all the network traffic. The key challenge in the Distinct Elements problem is the following: given a small set of samples where most of the colors are not observed, how to accurately extrapolate the number of unseens?

1.1 Main results

The fundamental limit of the Distinct Elements problem is characterized by the sample complexity, i.e., the smallest sample size needed to estimate the number of distinct colors with a prescribed accuracy and confidence level. A formal definition is the following:

Definition 1. The sample complexity $n^*(k, \Delta)$ is the minimal sample size n such that there exists an integer-valued estimator \hat{C} based on n balls drawn independently with replacements from the urn, such that $\mathbb{P}[|\hat{C} - C| \geq \Delta] \leq 0.1$ for any urn containing k balls with C different colors.¹

The main results of this paper provide bounds and constant-factor approximations of the sample complexity in various regimes summarized in Table 1, as well as computationally efficient algorithms. Below we highlight a few important conclusions drawn from Table 1:

From linear to sublinear samples: From the sample complexity for $k^{0.5+\delta} \leq \Delta \leq ck$ in Table 1, we conclude that the sample complexity is sublinear in k if and only if $\Delta = k^{1-o(1)}$, which also holds for sampling without replacement. To estimate within a constant fraction of balls $\Delta = ck$ for any small constant c , the sample complexity is $\Theta(\frac{k}{\log k})$, which coincides with the general support size estimation problem [VV11a, WY15] (see Section 1.2 for a detailed comparison). However, in other regimes we can achieve better performance by exploiting the discrete nature of the Distinct Elements problem.

From linear to superlinear samples: The transition from linear to superlinear sample complexity occurs near $\Delta = \sqrt{k}$. Although the exact sample complexity near $\Delta = \sqrt{k}$ is not completely resolved in the current paper, the lower bound and upper bound in Table 1 differ by a factor of at most $\log \log k$. In particular, the estimator via interpolation can achieve $\Delta = \sqrt{k}$ with $n = O(k \log \log k)$ samples, and achieving a precision of $\Delta \leq k^{0.5-o(1)}$ requires strictly superlinear sample size.

| Δ | Lower bound | Upper bound | Estimator |
|--|--|---|--|
| ≤ 1 | $\Theta(k \log k)$ | | Naïve |
| $[1, \sqrt{k}(\log k)^{-\delta}]$ | $\Theta(k \log \frac{k}{\Delta^2})$ | | Interpolation (Section 2.4) |
| $[\sqrt{k}(\log k)^{-\delta}, k^{0.5+\delta}]$ | $\Omega(k(1 \vee \log \frac{k}{\Delta^2}))$ | $O(k \log \frac{\log k}{1 \vee \log \frac{\Delta^2}{k}})$ | |
| $[k^{0.5+\delta}, ck]$ | $\Theta(\frac{k}{\log k} \log \frac{k}{\Delta})$ | | ℓ_2 -approximation (Section 2.2) |
| $[ck, (0.5 - \delta)k]$ | $k \exp(-\sqrt{O(\log k \log \log k)})^2$ | $O(\frac{k}{\log k})$ | |

Table 1: Summary of the sample complexity $n^*(k, \Delta)$, where δ is any sufficiently small constant and c is an absolute positive constant. The estimators are linear with coefficients obtained from either interpolation or ℓ_2 -approximation.

To establish the sample complexity, our lower bounds are obtained under zero-one loss and our upper bounds are under the (stronger) quadratic loss. Hence we also obtain the following

¹Clearly, since $\hat{C} - C \in \mathbb{Z}$, we shall assume without loss of generality that $\Delta \in \mathbb{N}$, with $\Delta = 1$ corresponding to the exact estimation of the number of distinct elements.

²A more precise result from [RRSS09] is the following: for $\Delta \in [ck, 0.5k - 2k^{3/4}\sqrt{\log k}]$, $n^*(k, \Delta) \geq k \exp(-\sqrt{O(\log k(\log \log k + \log \frac{k}{k/2-\Delta}))})$.

characterization of the minimax mean squared error (MSE) of the Distinct Elements problem:

$$\min_{\hat{C}} \max_{k\text{-ball urn}} \mathbb{E} \left(\frac{\hat{C} - C}{k} \right)^2 = \exp \left(-\Theta \left(1 \vee \frac{n \log k}{k} \wedge \log k \vee \frac{n}{k} \right) \right) \\ = \begin{cases} \Theta(1), & n \leq \frac{k}{\log k}, \\ \exp(-\Theta(\frac{n \log k}{k})), & \frac{k}{\log k} \leq n \leq k, \\ \exp(-\Theta(\log k)), & k \leq n \leq k \log k, \\ \exp(-\Theta(\frac{n}{k})), & n \geq k \log k, \end{cases}$$

where \hat{C} denotes an estimator using n samples with replacements and C is the number of distinct colors in a k -ball urn.

1.2 Related work

Statistics literature The Distinct Elements problem is equivalent to estimating the number of species (or classes) in a finite population, which has been extensively studied in the statistics (see surveys [BF93, GS04]) and the numismatics literature (see survey [Est86]). Motivated by various practical applications, a number of statistical models have been introduced for this problem, the most popular four being (cf. [BF93, Figure 1]):

- *Multinomial model*: n samples are drawn uniformly at random with replacement;
- *Hypergeometric model*: n samples are drawn uniformly at random without replacement;
- *Bernoulli model*: each individual is observed independently with some fixed probability, and thus the total number of samples is a binomial random variable;
- *Poisson model*: the number of observed samples in each class is independent and Poisson distributed, and thus the total sample size is also a Poisson random variable.

These models are closely related: conditioned on the sample size, the Bernoulli model coincides with the hypergeometric one, and Poisson model coincides with the multinomial one; furthermore, hypergeometric model can simulate multinomial one and is hence more informative. The multinomial model is adopted as the main focus of this paper and the sample complexity in Definition 1 refers to the number of samples with replacement. In the undersampling regime where the sample size is significantly smaller than the population size, all four models are approximately equivalent. See Appendix A for a rigorous justification and detailed comparisons.

Under these models various estimators have been proposed such as unbiased estimators [Goo49], Bayesian estimators [Hil79], variants of Good-Turing estimators [CL92], etc. None of these methodologies, however, have a provable worst-case guarantee. Finally, we mention a closely related problem of estimating the number of connected components in a graph based on sampled induced subgraphs. In the special case where the underlying graph consists of disjoint cliques, the problem is exactly equivalent to the Distinct Elements problem [Fra78].

Computer science literature The interests in the Distinct Elements problem also arise in the database literature, where various intuitive estimators [HOT88, NS90] have been proposed under simplifying assumptions such as uniformity, and few performance guarantees are available. More recent work in [CCMN00, BYKS01] obtained the optimal sample complexity under the *multiplicative*

error criterion, where the minimum sample size to estimate the number of distinct elements within a factor of $1 \pm \epsilon$ is shown to be $\Theta(k/\epsilon^2)$. For this task, it turns out the least favorable scenario is to distinguish an urn with unitary color from one with *almost* unitary color, the impossibility of which implies large multiplicative error. However, the optimal estimator performs poorly compared with others on an urn with many distinct colors [CCMN00], the case where most estimators enjoy small multiplicative error. In view of the limitation of multiplicative error, additive error is later considered by [RRSS09, Val11]. To achieve an additive error of ck for a constant $c \in (0, \frac{1}{2})$, the result in [CCMN00] only implies an $\Omega(1/c)$ sample complexity lower bound, whereas a much stronger lower bound scales $k^{1-O(\sqrt{\frac{\log \log k}{\log k}})}$ is obtained in [RRSS09] that scales almost linearly. Determining the optimal sample complexity under additive error is the focus of the present paper.

The Distinct Elements problem can be viewed as a special case of the Support Size problem, where the goal is to estimate the cardinality of the support of an unknown discrete distribution, whose nonzero probabilities are at least $\frac{1}{k}$, based on independent samples. Improving previous results in [VV11a], the optimal sample complexity has been recently determined in [WY15] to be

$$\Theta\left(\frac{k}{\log k} \log^2 \frac{k}{\Delta}\right). \quad (1)$$

Samples drawn from a k -ball urn with replacement can be viewed as i.i.d. samples from a distribution supported on the set $\{\frac{1}{k}, \frac{2}{k}, \dots, \frac{k}{k}\}$. From this perspective, any support size estimator, as well as its performance guarantee, is applicable to the Distinct Elements problem.

We briefly describe and compare the strategy to construct estimators in [WY15] and the current paper. Both are based on the idea of *polynomial approximation*, a powerful tool to circumvent the nonexistence of unbiased estimators [LNS99]. The key is to approximate the function to be estimated by a polynomial, whose degree is chosen to balance the approximation error (bias) and the estimation error (variance). The worst-case performance guarantee for the Support Size problem in [WY15] is governed by the uniform approximation error over an interval where the probabilities may reside. In contrast, in the Distinct Elements problem, samples are generated from a distribution supported on a *discrete* set of values. Uniform approximation over a discrete subset leads to smaller approximation error and, in turn, improved sample complexity. It turns out that $O(\frac{k}{\log k} \log \frac{k}{\Delta})$ samples are sufficient to achieve an additive error of Δ that satisfies $k^{0.5+O(1)} \leq \Delta \leq O(k)$, which strictly improves the sample complexity (1) for the Support Size problem, thanks to the discrete structure of the Distinct Elements problem.

The Distinct Elements problem considered here is not to be confused with the formulation in the streaming literature, where the goal is to approximate the number of distinct elements in the observations with low space complexity, see, e.g., [FFGM07, KNW10]. Their algorithms optimize the memory consumption, but still require a full pass of every ball in the urn. This is different from the setting in the current paper, where only random samples drawn from the urn are available.

1.3 Organization

The paper is organized as follows: In Section 2 we describe a unified approach to construct estimators via discrete polynomial approximation, whose bias is analyzed in Section 2.2 and variance is upper bounded in Sections 2.3 and 2.4 separately. In Section 3 we obtain lower bounds on the sample complexity in Table 1 which establish the optimality of the proposed estimators. Connections between the four sampling model mentioned in Section 1.2 is detailed in Appendix A. Proof of auxiliary results are deferred to Appendix B and Appendix C. Finally, Appendix D explains how sample complexity bounds summarized in Table 1 follow from results in Sections 2 and 3.

1.4 Notations

All logarithms are with respect to the natural base. The transpose of a matrix A is denoted by A' . Let $\mathbf{1}$ denote the all-one column vector. Let $\|\cdot\|_p$ denote the vector ℓ_p -norm, for $1 \leq p \leq \infty$. Let $\text{Poi}(\lambda)$ be the Poisson distribution with mean λ , $\text{Bern}(p)$ be the Bernoulli distribution with mean p , $\text{Binomial}(n, p)$ be the binomial distribution with n trials and success probability p , and $\text{Hypergeometric}(N, K, n)$ be the hypergeometric distribution with probability mass function $\binom{K}{k} \binom{N-K}{n-k} / \binom{N}{n}$, for $0 \vee (n+K-N) \leq k \leq n \wedge K$. The n -fold product of a distribution P is denoted by $P^{\otimes n}$. We use standard big- O notations: for any positive sequence $\{a_n\}$ and $\{b_n\}$, $a_n = O(b_n)$ or $a_n \lesssim b_n$ if $a_n \leq Cb_n$ for some absolute constant $C > 0$; $a_n = \Omega(b_n)$ or $a_n \gtrsim b_n$ if $b_n = O(a_n)$; $a_n = \Theta(b_n)$ or $a_n \asymp b_n$ if both $a_n = O(b_n)$ and $b_n = O(a_n)$; $a_n = o(b_n)$ if $\lim a_n/b_n = 0$; $a_n = \omega(b_n)$ if $b_n = o(a_n)$. Furthermore, we use $o_n(1)$ to indicate convergence in n that is uniform in all other parameters.

2 Linear estimators via discrete polynomial approximation

In this section we develop a unified framework to construct linear estimators and analyze its performance. Note that linear estimators have been previously used for estimating distribution functionals [Pan04, VV11a, VV11b, WY15]. As commonly done in the literature, we assume the *Poisson sampling model*, where the sample size is a random variable $\text{Poi}(n)$ instead of being exactly n . Furthermore, the histograms of the samples are independent which simplifies the analysis. Any estimator under the Poisson sampling model can be easily modified for fixed sample size, and vice versa, thanks to the concentration of the Poisson random variable near its mean. Consequently, the sample complexities of these two models are close to each other, as shown in Corollary 1 in Appendix A.

2.1 Guarantees on general linear estimators

Recall that C denotes the number of distinct colors in a urn containing k colored balls. Let k_i denote the number of balls of the i th color in the urn. Then $\sum_i k_i = k$ and $C = \sum_i \mathbf{1}_{\{k_i > 0\}}$. Let X_1, X_2, \dots be independently drawn with replacements from the urn. Equivalently, X_i 's are iid distributed according to a distribution $P = (p_i)$, where $p_i = k_i/k$ is the fraction of balls of the i th color. The observed data are X_1, \dots, X_N , where the sample size N is independent and distributed as $\text{Poi}(n)$. Under the Poisson model (or any of the sampling models described in Section 1.2), the histograms $\{N_i\}$ are sufficient statistics for inferring any aspect of the urn configuration. Here N_i is the number of balls of the i th color observed in the sample, which are independently distributed as $\text{Poi}(np_i)$. Furthermore, the *fingerprints* $\{\Phi_j\}_{j \geq 1}$, which are the histogram of the histograms, are also sufficient for estimating any permutation-invariant distributional property [Pan03], in particular, the number of colors. Specifically, the j th fingerprint Φ_j denotes the number of colors that appear exactly j times. Note that $U \triangleq \Phi_0$, the number of unseen colors, is not observed.

The naïve estimator, “what you see is what you get”, is simply the number of observed distinct colors, which can be expressed in terms of fingerprints as

$$\hat{C}_{\text{seen}} = \sum_{j \geq 1} \Phi_j,$$

which is typically an underestimate since $C = \hat{C}_{\text{seen}} + U$. In turn, our estimator is

$$\tilde{C} = \hat{C}_{\text{seen}} + \hat{U}, \tag{2}$$

which adds a linear correction term

$$\hat{U} = \sum_{j \geq 1} u_j \Phi_j, \quad (3)$$

where the coefficients u_j 's are to be specified. Since the fingerprints Φ_0, Φ_1, \dots are dependent (for example, they sum up to C), (3) serves as a linear predictor of $U = \Phi_0$ in terms of the observed fingerprints. Equivalently, in terms of histograms, the estimator has the following decomposable form:

$$\tilde{C} = \sum_i g(N_i), \quad (4)$$

where $g : \mathbb{Z}_+ \rightarrow \mathbb{R}$ satisfies $g(0) = 0$ and $g(j) = 1 + u_j$ for $j \geq 1$.

The main idea to choose the coefficients u_j is to achieve a good trade-off between the variance and the bias. Note that linear estimators can easily achieve exactly zero bias, which, however, comes at the price of high variance. To see this, note that the bias of the estimator (4) is $\mathbb{E}[\tilde{C}] - C = \sum \mathbb{E}[g(N_i)] - 1$, where

$$|\mathbb{E}[g(N_i) - 1]| = e^{-np_i} \left| -1 + \sum_{j=1}^L k_i^j \frac{u_j (n/k)^j}{j!} \right| \leq e^{-n/k} \max_{a \in [k]} |\phi(a) - 1|. \quad (5)$$

where $\phi(a) \triangleq \sum_{j \geq 1} a^j \frac{u_j (n/k)^j}{j!}$ is a (formal) power series with $\phi(0) = 0$. The right-hand side of (5) can be made zero by choosing ϕ to be, e.g., the Lagrange interpolating polynomial that satisfies $\phi(0) = -1$ and $\phi(i) = 0$ for $i \in [k]$; however, this strategy results in a high-degree polynomial ϕ with large coefficients, which, in turn, lead to large variance of the estimator.

To reduce the variance of our estimator, we only use the first L fingerprints in (3) by setting $u_j = 0$ for all $j > L$, where L is chosen to be proportional to $\log k$. This restricts the polynomial degree in (5) to at most L and, while possibly incurring bias, reduces the variance. Another reason for only using the first few fingerprints is that higher-order fingerprints are almost uncorrelated with the number of unseens Φ_0 . For instance, if red balls are observed for $n/2$ times, the only information this reveals is that approximately half of the urn are red. In fact, the correlation between Φ_0 and Φ_j decays exponentially (see Appendix B). Therefore for $L = \Theta(\log k)$, $\{\Phi_j\}_{j > L}$ offer little predictive power about Φ_0 . Moreover, if a color is observed at most L times, say, $N_i \leq L$, this implies that, with high probability, $k_i \leq M$, where $M = O(kL/n)$, thanks to the concentration of Poisson random variables. Therefore, effectively we only need to consider those colors that appear in the urn for at most M times, i.e., $k_i \in [M]$, for which the bias is at most

$$|\mathbb{E}[g(N_i) - 1]| \leq e^{-n/k} \max_{a \in [M]} |\phi(a) - 1| = e^{-n/k} \max_{x \in [M]/M} |p(x) - 1| = e^{-n/k} \|Bw - \mathbf{1}\|_\infty, \quad (6)$$

where $p(x) \triangleq \phi(Mx) = \sum_{j=1}^L w_j x^j$, $w = (w_1, \dots, w_L)'$, and

$$w_j \triangleq \frac{u_j (Mn/k)^j}{j!}, \quad B \triangleq \begin{pmatrix} 1/M & (1/M)^2 & \dots & (1/M)^L \\ 2/M & (2/M)^2 & \dots & (2/M)^L \\ \vdots & \vdots & \ddots & \vdots \\ 1 & 1 & \dots & 1 \end{pmatrix} \quad (7)$$

is a (partial) Vandermonde matrix. Lastly, since $\hat{C}_{\text{seen}} \leq C \leq k$, we define the final estimator to be \tilde{C} projected to the interval $[\hat{C}_{\text{seen}}, k]$. We have the following error bound:

Proposition 1. Assume the Poisson sampling model. Let

$$L = \alpha \log k, \quad M = \frac{\beta k \log k}{n}, \quad (8)$$

for any $\beta > \alpha$ such that L and M are integers. Let $w \in \mathbb{R}^L$. Let \tilde{C} be defined in (2) with $u_j = w_j j! (\frac{k}{nM})^j$ for $j \in [L]$ and $u_j = 0$ otherwise. Define $\hat{C} \triangleq (\tilde{C} \vee \hat{C}_{\text{seen}}) \wedge k$. Then

$$\mathbb{E}(\hat{C} - C)^2 \leq k^2 e^{-2n/k} \|Bw - \mathbf{1}\|_\infty^2 + k e^{-n/k} + k \max_{m \in [M]} \mathbb{E}_{N \sim \text{Poi}(nm/k)}[u_N^2] + k^{-(\beta - \alpha \log \frac{e\beta}{\alpha} - 3)}. \quad (9)$$

Proof. Since $\hat{C}_{\text{seen}} \leq C \leq k$, \hat{C} is always an improvement of \tilde{C} . Define the event $E = \cap_{i=1}^k \{N_i \leq L \Rightarrow kp_i \leq M\}$. Since $\beta > \alpha$, applying the Chernoff bound and the union bound yields $\mathbb{P}[E^c] \leq k^{1-\beta+\alpha \log \frac{e\beta}{\alpha}}$, and thus

$$\mathbb{E}(\hat{C} - C)^2 \leq \mathbb{E}((\hat{C} - C)\mathbf{1}_E)^2 + k^2 \mathbb{P}[E^c] \leq \mathbb{E}((\tilde{C} - C)\mathbf{1}_E)^2 + k^{3-\beta+\alpha \log \frac{e\beta}{\alpha}}. \quad (10)$$

The decomposable form of \tilde{C} in (4) leads to

$$(\tilde{C} - C)\mathbf{1}_E = \sum_{i:k_i \in [M]} (g(N_i) - 1)\mathbf{1}_{\{N_i \leq L\}} \triangleq \mathcal{E}.$$

In view of the bias analysis in (6), we have

$$|\mathbb{E}[\mathcal{E}]| \leq \sum_{i:k_i \in [M]} e^{-nk_i/k} \|Bw - \mathbf{1}\|_\infty \leq k e^{-n/k} \|Bw - \mathbf{1}\|_\infty. \quad (11)$$

Recall that $g(0) = 0$ and $g(j) = u_j + 1$ for $j \in [L]$. Since $N_i \stackrel{\text{ind}}{\sim} \text{Poi}(nk_i/k)$, we have

$$\begin{aligned} \text{var}[\mathcal{E}] &= \sum_{i:k_i \in [M]} \text{var}[(g(N_i) - 1)\mathbf{1}_{\{N_i \leq L\}}] \leq \sum_{i:k_i \in [M]} \mathbb{E}[(g(N_i) - 1)^2 \mathbf{1}_{\{N_i \leq L\}}] \\ &= \sum_{i:k_i \in [M]} \left(e^{-nk_i/k} + \mathbb{E}[u_{N_i}^2] \right) \leq k e^{-n/k} + k \max_{m \in [M]} \mathbb{E}_{N \sim \text{Poi}(nm/k)}[u_N^2]. \end{aligned} \quad (12)$$

Combining the upper bound on the bias in (11) and the variance in (12) yields an upper bound on $\mathbb{E}[\mathcal{E}^2]$. Then the MSE in (9) follows from (10). \square

Proposition 1 suggests that the coefficients of the linear estimator can be chosen by solving the following linear programming (LP)

$$\min_{w \in \mathbb{R}^L} \|Bw - \mathbf{1}\|_\infty \quad (13)$$

and showing that the solution does not have large entries. Instead of the ℓ_∞ -approximation problem (13), whose optimal value is difficult to analyze, we solve the ℓ_2 -approximation problem as a relaxation:

$$\min_{w \in \mathbb{R}^L} \|Bw - \mathbf{1}\|_2, \quad (14)$$

which is an upper bound of (13), and is in fact within an $O(\log k)$ factor since $M = O(k \log k/n)$ and $n = \Omega(k/\log k)$. In the remainder of this section, we consider two separate cases:

- $M > L$ ($n \lesssim k$): In this case, the linear system in (14) is overdetermined and the minimum is non-zero. Surprisingly, as shown in Section 2.2, the exact optimal value can be found in close form using discrete orthogonal polynomials. The coefficients of the solution can be bounded using the minimum singular value of the matrix B , which is analyzed in Section 2.3.

- $M \leq L$ ($n \gtrsim k$): In this case, the linear system is underdetermined and the minimum in (14) is zero. To bound the variance, it turns out that the coefficients bound obtained from the minimum singular value is not precise enough in the regime. Instead, we express the coefficients in terms of Lagrange interpolating polynomials and use Stirling numbers to obtain sharp variance bounds. This analysis is carried out in Section 2.4.

We finish this subsection with two remarks:

Remark 1 (Discrete versus continuous approximation). The optimal estimator for the **Support Size** problem in [WY15] has the same linear form as (2); however, since the probabilities can take any values in an interval, the coefficients are found to be the solution of the continuous polynomial approximation problem

$$\inf_p \max_{x \in [\frac{1}{M}, 1]} |p(x) - 1| = \exp\left(-\Theta\left(\frac{L}{\sqrt{M}}\right)\right). \quad (15)$$

where the infimum is taken over all degree- L polynomials such that $p(0) = 0$, achieved by the (appropriately shifted and scaled) Chebyshev polynomial [Tim63]. In contrast, in Section 2.2 we show that the following discrete version, which is equivalent to the LP (13), satisfies

$$\inf_p \max_{x \in \{\frac{1}{M}, \frac{2}{M}, \dots, 1\}} |p(x) - 1| = \text{poly}(M) \exp\left(-\Theta\left(\frac{L^2}{M}\right)\right), \quad (16)$$

provided $L < M$. The difference between (15) and (16) explains why the sample complexity (1) for the **Support Size** problem has an extra log factor compared to that of the **Distinct Elements** problem in Table 1. When the sample size n is large enough, interpolation is used in lieu of approximation. See Fig. 1 for an illustration.

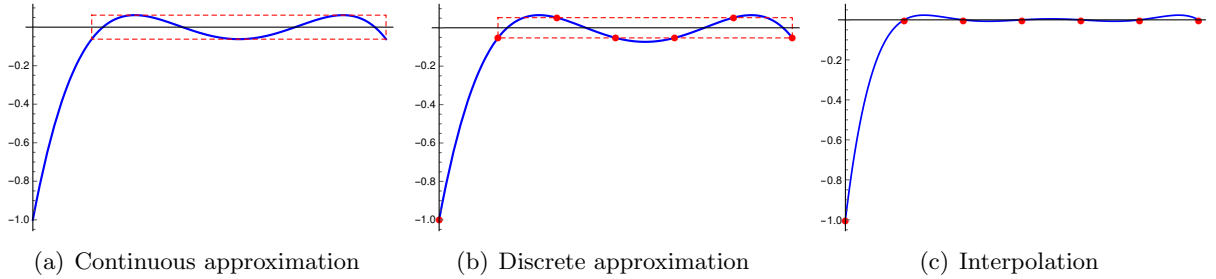


Figure 1: Continuous and discrete polynomial approximations for $M = 6$ and degree $L = 4$, where (a) and (b) shows optimal solution to (15) and (16) respectively. The interpolating polynomial in (c) requires a higher degree $L = 6$.

Remark 2 (Time complexity). The time complexity of the estimator (2) consists of: (a) Computing histograms N_i and fingerprints Φ_j of n samples: $O(n)$; (b) Computing the coefficients w by solving the least square problem in (6): $O(L^2(M+L))$; (c) Evaluating the linear combination (2): $O(n \wedge k)$. As shown in Table 1, for accurate estimation the sample complexity is $n = \Omega(\frac{k}{\log k})$, which implies $L = O(\log k)$ and $M = O(\log^2 k)$. Therefore, the overall time complexity is $O(n + \log^4 k)$:

2.2 Exact solution to the ℓ_2 -approximation

Next we give an explicit solution to the ℓ_2 -approximation problem (14). In general, the optimal solution w^* is given by $w^* = (B'B)^{-1}B'\mathbf{1}$ and the minimum value is the Euclidean distance between

the all-one vector $\mathbf{1}$ and the column span of B , which, in the case of $M > L$, is non-zero. Taking advantage of the Vandermonde structure of the matrix B in (7), we note that (14) can be interpreted as finding the orthogonal projection of the constant function onto the linear space of polynomials of degree between 1 and L defined on the discrete set $[M]/M$. Using the orthogonal polynomials with respect to the counting measure, known as *discrete Chebyshev (or Gram) polynomials* (see [Sze75, Section 2.8] or [NUS91, Section 2.4.2]), we show that, surprisingly, the optimal value of the ℓ_2 -approximation can be found in close form:

Lemma 1. *For all $L \geq 1$ and $M \geq L + 1$,*

$$\min_{w \in \mathbb{R}^L} \|Bw - \mathbf{1}\|_2 = \left[\frac{\binom{M+L+1}{L+1}}{\binom{M}{L+1}} - 1 \right]^{-1/2} = \left[\exp \left(\Theta \left(\frac{L^2}{M} \right) \right) - 1 \right]^{-1/2}. \quad (17)$$

Proof. Define the following inner product between functions f and g :

$$\langle f, g \rangle \triangleq \sum_{i=1}^M f \left(\frac{i}{M} \right) g \left(\frac{i}{M} \right) \quad (18)$$

and the induced norm $\|f\| \triangleq \sqrt{\langle f, f \rangle}$. The least square problem (17) can be equivalently formulated as

$$\min_w \|-1 + w_1x + w_2x^2 + \dots + w_Lx^L\|. \quad (19)$$

This can be analyzed using the orthogonal polynomials under the inner product (18), which we describe next.

Recall the discrete Chebyshev polynomial [Sze75, Sec. 2.8]:

$$t_m(x) \triangleq \frac{1}{m!} \Delta^m p_m(x) = \frac{1}{m!} \sum_{j=0}^m (-1)^j \binom{m}{j} p_m(x + m - j), \quad 0 \leq m \leq M-1, \quad (20)$$

where

$$p_m(x) \triangleq x(x-1) \cdots (x-m+1)(x-M)(x-M-1) \cdots (x-M-m+1), \quad (21)$$

and Δ^m denotes the m -th order forward difference. The polynomials $\{t_0, \dots, t_{M-1}\}$ are orthogonal with respect to the counting measure over the discrete set $\{0, 1, \dots, M-1\}$; in particular, we have (cf. [Sze75, Sec. 2.8.2, 2.8.3]):

$$\begin{aligned} \sum_{x=0}^{M-1} t_m(x) t_l(x) &= 0, \quad m \neq l, \\ \sum_{x=0}^{M-1} t_m^2(x) &= c(M, m) \triangleq \frac{M(M^2 - 1^2)(M^2 - 2^2) \cdots (M^2 - m^2)}{2m + 1}. \end{aligned}$$

By appropriately shifting and scaling the set of polynomials t_m , we define an orthonormal basis for the set of polynomials of degree at most $L \leq M-1$ under the inner product (18) by

$$\phi_m(x) = \frac{t_m(Mx - 1)}{\sqrt{c(M, m)}}, \quad m = 0, \dots, L. \quad (22)$$

Using this basis, the least square problem (19) can be equivalently formulated as

$$\min_{a: \sum_{i=1}^L a_i \phi_i(0) = -1} \left\| \sum_{i=0}^L a_i \phi_i \right\| = \min_{a: a' \phi(0) = -1} \|a\|_2,$$

where $\phi(0) \triangleq (\phi_0(0), \dots, \phi_L(0))$ and $a = (a_0, \dots, a_L)$. The optimal value is clearly $\frac{1}{\|\phi(0)\|_2}$, achieved by $a^* = -\frac{\phi(0)}{\|\phi(0)\|_2^2}$.

From (21) we have $p_m(0) = p_m(1) = \dots = p_m(m-1) = 0$. By the formula of t_m in (20), we obtain

$$t_m(-1) = \frac{1}{m!}(-1)^m p_m(-1) = (-1)^m \prod_{j=1}^m (M+j).$$

In view of the definition of ϕ_m in (22), we have

$$\phi_m(0) = \frac{t_m(-1)}{\sqrt{c(M, m)}} = \frac{(-1)^m \prod_{j=1}^m (M+j)}{\sqrt{\frac{M \prod_{j=1}^m (M^2-j^2)}{2m+1}}} = (-1)^m \sqrt{\frac{2m+1}{M} \prod_{j=1}^m \frac{M+j}{M-j}}.$$

Therefore

$$\|\phi(0)\|_2^2 = \sum_{m=0}^L \frac{2m+1}{M} \prod_{j=1}^m \frac{M+j}{M-j} = \frac{\binom{M+L+1}{L+1}}{\binom{M}{L+1}} - 1,$$

where the last equality follows from induction since

$$\frac{\binom{M+L+1}{L+1}}{\binom{M}{L+1}} - \frac{\binom{M+L}{L}}{\binom{M}{L}} = \frac{2L+1}{M} \prod_{j=1}^L \frac{M+j}{M-j}.$$

This proves the first equality in (17).

The second equality in (17) is a direct consequence of Stirling's approximation. If $M = L+1$, then

$$\frac{\binom{M+L+1}{L+1}}{\binom{M}{L+1}} = \binom{2(L+1)}{L+1} = \exp(\Theta(L)). \quad (23)$$

If $M \geq L+2$, denoting $x = \frac{L+1}{M}$ and applying $n! = \sqrt{2\pi n} \left(\frac{n}{e}\right)^n (1 + \Theta(\frac{1}{n}))$ when $n \geq 1$, we have

$$\begin{aligned} \frac{\binom{M+L+1}{L+1}}{\binom{M}{L+1}} &= \frac{(M+L+1)!(M-L-1)!}{(M!)^2} = \frac{(M(1+x))!(M(1-x))!}{(M!)^2} \\ &= \frac{\sqrt{2\pi M(1+x)} \left(\frac{M(1+x)}{e}\right)^{M(1+x)} \sqrt{2\pi M(1-x)} \left(\frac{M(1-x)}{e}\right)^{M(1-x)} (1 + \Theta(\frac{1}{M(1+x)} + \frac{1}{M(1-x)}))}{2\pi M \left(\frac{M}{e}\right)^{2M} (1 + \Theta(\frac{1}{M}))} \\ &= \sqrt{1-x^2} \exp(M((1+x)\log(1+x) + (1-x)\log(1-x))) \frac{1 + \Theta(\frac{1}{M(1-x^2)})}{1 + \Theta(\frac{1}{M})} \\ &= \exp\left(\Theta(Mx^2) + \frac{1}{2}\log(1-x^2) + \log \frac{1 + \Theta(\frac{1}{M(1-x^2)})}{1 + \Theta(\frac{1}{M})}\right), \end{aligned} \quad (24)$$

where the last step follows from $(1+x)\log(1+x) + (1-x)\log(1-x) = \Theta(x^2)$ when $0 \leq x \leq 1$. In the exponent of (24), the term $\Theta(Mx^2)$ dominates when $M \geq L+2$. Applying (23) and (24) to the exact solution (17) yields the desired approximation. \square

2.3 Minimum singular values of real rectangle Vandermonde matrices

In Proposition 1 the variance of our estimator is bounded by the magnitude of coefficients u , which is related to the polynomial coefficients w by (7). A classical result from approximation theory is that if a polynomial is bounded over a compact interval, its coefficients are at most exponential in the degree [Tim63, 2.9.11]: for any degree- L polynomial $p(x) = \sum_{i=0}^L w_i x^i$,

$$\max_{0 \leq i \leq L} |w_i| \leq \max_{x \in [0,1]} |p(x)| \exp(O(L)), \quad (25)$$

which is tight when p is the Chebyshev polynomial. This fact has been applied in statistical contexts to control the variance of estimators obtained from best polynomial approximation [CL11, WY16, WY15]. In contrast, for the **Distinct Elements** problem, the polynomial is only known to be bounded over the discretized interval. Nevertheless, we show that the bound (25) continues to hold as long as the discretization level exceeds the degree:

$$\max_{0 \leq i \leq L} |w_i| \leq \max_{x \in \{\frac{1}{M}, \frac{2}{M}, \dots, 1\}} |p(x)| \exp(O(L)), \quad (26)$$

provided that $M \geq L + 1$ (see Remark 3 after Lemma 2). Clearly, (26) implies (25) by sending $M \rightarrow \infty$. If $M \leq L$, coefficient bound like (26) is impossible, because one can add to p an arbitrary degree- L interpolating polynomial that evaluates to zero at all M points.

To bound the coefficients, note that the optimal solution of ℓ_2 -approximation is $w^* = (B'B)^{-1}B'\mathbf{1}$, and consequently

$$\|w^*\|_2 \leq \frac{\|\mathbf{1}\|_2}{\sigma_{\min}(B)}, \quad (27)$$

where $\sigma_{\min}(B)$ denotes the smallest singular value of B . Let

$$\bar{B} \triangleq [\mathbf{1}, B] = \begin{pmatrix} 1 & 1/M & (1/M)^2 & \dots & (1/M)^L \\ 1 & 2/M & (2/M)^2 & \dots & (2/M)^L \\ 1 & \vdots & \vdots & \ddots & \vdots \\ 1 & 1 & 1 & \dots & 1 \end{pmatrix}$$

which is an $M \times (L+1)$ Vandermonde matrix and satisfies $\sigma_{\min}(\bar{B}) \leq \sigma_{\min}(B)$. The Gram matrix of \bar{B} is an instance of *moment matrices*. A moment matrix associated with a probability measure μ is a Hankel matrix M given by $M_{i,j} = m_{i+j-2}$, where $m_\ell = \int x^\ell d\mu$ denotes the ℓ th moment of μ . Then $\frac{1}{M}\bar{B}'\bar{B}$ is the moment matrix associated with the uniform distribution over the discrete set $\{\frac{1}{M}, \frac{2}{M}, \dots, 1\}$, which converges to the uniform distribution over the interval $(0,1)$. The moment matrix of the uniform distribution is the famous *Hilbert matrix* H , with

$$H_{ij} = \frac{1}{i+j-1}$$

which is a well-studied example of ill-conditioned matrices in the numerical analysis literature. In particular, it is known that the condition number of the $L \times L$ Hilbert matrix is $O(\frac{(1+\sqrt{2})^{4L}}{\sqrt{L}})$ [Tod54] and the operator norm is $\Theta(1)$, and thus the minimum singular value is exponentially small in the degree. Therefore we expect the discrete moment matrix $\frac{1}{M}\bar{B}'\bar{B}$ behaves similarly to the Hilbert matrix when M is large enough. Interestingly, we show that this is indeed the case as soon as M exceeds L (otherwise the minimum singular value is zero).

Lemma 2. For all $M \geq L + 1$,

$$\sigma_{\min} \left(\frac{\bar{B}}{\sqrt{M}} \right) \geq \frac{1}{L^2 2^{7L} (2L + 1)} \left(\frac{M + L}{eM} \right)^{L+0.5}. \quad (28)$$

Remark 3. The inequality (26) follows from Lemma 2 since the coefficient vector $w = (w_0, \dots, w_L)$ satisfies $\|w\|_\infty \leq \|w\|_2 \leq \frac{1}{\sigma_{\min}(\bar{B})} \|\bar{B}w\|_2 \leq \frac{\sqrt{M}}{\sigma_{\min}(\bar{B})} \|\bar{B}w\|_\infty$.

Remark 4. The extreme singular values of square Vandermonde matrices have been extensively studied (c.f. [Gau90, Bec00] and the references therein). For rectangular Vandermonde matrices, the focus was mainly with nodes on the unit circle in the complex domain [CGR90, Fer99, Moi15] with applications in signal processing. In contrast, Lemma 2 is on rectangle Vandermonde matrices with real nodes. The result on integers nodes in [EPS01] turns out to be too crude for the purpose of this paper.

Proof. Note that $\bar{B}'\bar{B}$ is the Gramian of monomials $\mathbf{x} = (1, x, x^2, \dots, x^L)'$ under the inner product defined in (18). When $M \geq L + 1$, the orthonormal basis $\phi = (\phi_0, \dots, \phi_L)'$ under the inner product (18) are given in (22). Let $\phi = \mathbf{L}\mathbf{x}$ where $\mathbf{L} \in \mathbb{R}^{(L+1) \times (L+1)}$ is a lower triangular matrix and \mathbf{L} consists of the coefficients of ϕ . Taking the Gramian of ϕ yields that $I = \mathbf{L}(\bar{B}'\bar{B})\mathbf{L}'$, i.e., \mathbf{L}^{-1} can be obtained from the Cholesky decomposition: $\bar{B}'\bar{B} = (\mathbf{L}^{-1})(\mathbf{L}^{-1})'$. Then³

$$\sigma_{\min}^2(\bar{B}) = \frac{1}{\|\mathbf{L}\|_{op}^2} \geq \frac{1}{\|\mathbf{L}\|_F^2}, \quad (29)$$

where $\|\cdot\|_{op}$ denotes the ℓ_2 operator norm, which is the largest singular value of L , and $\|\cdot\|_F$ denotes the Frobenius norm. By definition, $\|\mathbf{L}\|_F^2$ is the sum of all squared coefficients of ϕ_0, \dots, ϕ_L . A useful method to bound the sum-of-squares of the coefficients of a polynomial is by its maximal modulus over the unit circle on the complex plane. Specifically, for any polynomial $p(z) = \sum_{i=0}^n a_i z^i$, we have

$$\sum_{i=0}^n |a_i|^2 = \frac{1}{2\pi} \oint_{|z|=1} |p(z)|^2 dz \leq \sup_{|z|=1} |p(z)|^2. \quad (30)$$

Therefore

$$\sigma_{\min}(\bar{B}) \geq \frac{1}{\|\mathbf{L}\|_F} \geq \frac{1}{\sqrt{\sum_{m=0}^L \sup_{|z|=1} |\phi_m(z)|^2}} \geq \frac{1}{\sqrt{L+1}} \frac{1}{\sup_{0 \leq m \leq L, |z|=1} |\phi_m(z)|}. \quad (31)$$

For a given M , the orthonormal basis $\phi_m(x)$ in (22) is proportional to the discrete Chebyshev polynomials $t_m(Mx - 1)$. The classical asymptotic result for the discrete Chebyshev polynomials shows that [Sze75, (2.8.6)]

$$\lim_{M \rightarrow \infty} M^{-m} t_m(Mx) = P_m(2x - 1),$$

where P_m is the Legendre polynomial of degree m . This gives the intuition that $t_m(x) \approx M^m$ for real-valued $x \in [0, M]$. We have the following non-asymptotic upper bound (proved in Appendix C) for t_m over the complex plane:

Lemma 3. For all $0 \leq m \leq M - 1$,

$$|t_m(z)| \leq m^2 2^{6m} \sup_{0 \leq \xi \leq m} (|z + \xi| \vee M)^m. \quad (32)$$

³The lower bound (29), which was also obtained in [CL99, (1.13)] using Cauchy-Schwartz inequality, is tight up to polynomial terms in view of the fact that $\|\mathbf{L}\|_F \leq (L+1)\|\mathbf{L}\|_{op}$

Applying (32) on the definition of ϕ_m in (22), for any $|z| = 1$ and any $M \geq L + 1$, we have

$$|\phi_m(z)| = \frac{|t_m(Mz - 1)|}{\sqrt{c(M, m)}} \leq \frac{m^2 2^{7m} M^m}{\sqrt{\frac{M(M^2-1^2)(M^2-2^2)\dots(M^2-m^2)}{2m+1}}}.$$

The right-hand side is increasing with m . Therefore,

$$\begin{aligned} \sup_{0 \leq m \leq L, |z|=1} |\phi_m(z)| &\leq \frac{L^2 2^{7L} M^L}{\sqrt{\frac{M(M^2-1^2)(M^2-2^2)\dots(M^2-L^2)}{2L+1}}} \\ &= \frac{1}{\sqrt{M}} L^2 2^{7L} \sqrt{2L+1} \sqrt{\frac{M^{2L+1}}{\binom{M+L}{2L+1} (2L+1)!}}. \end{aligned}$$

Combining (31), we obtain

$$\begin{aligned} \sigma_{\min} \left(\frac{\bar{B}}{\sqrt{M}} \right) &\geq \frac{1}{L^2 2^{7L} \sqrt{(L+1)(2L+1)}} \sqrt{\frac{\binom{M+L}{2L+1} (2L+1)!}{M^{2L+1}}} \\ &\geq \frac{1}{L^2 2^{7L} (2L+1)} \left(\frac{M+L}{eM} \right)^{L+0.5}, \end{aligned}$$

where in the last inequality we used $\binom{n}{k} \geq \left(\frac{n}{k}\right)^k$ and $n! \geq \left(\frac{n}{e}\right)^n$. \square

Using the optimal solution w^* to the ℓ_2 -approximation problem (14) as the coefficient of the linear estimator \hat{C} , the following performance guarantee is obtained by applying Lemma 1 and Lemma 2 to bound the bias and variance, respectively:

Theorem 1. *Assume the Poisson sampling model.*

$$\mathbb{E}(\hat{C} - C)^2 \leq k^2 \exp \left(-\Theta \left(1 \vee \frac{n \log k}{k} \wedge \log k \right) \right). \quad (33)$$

Proof. If $n \leq \frac{k}{\log k}$, then the upper bound in (33) is $\Theta(k^2)$, which is trivial thanks to the thresholds that $\hat{C} = (\tilde{C} \vee \hat{C}_{\text{seen}}) \wedge k$. It is hereinafter assumed that $n \geq \frac{k}{\log k}$, or equivalently $M \leq \frac{\beta}{\alpha^2} L^2$. Let α, β be constants to be specified. Then, from Lemma 1,

$$\|Bw^* - \mathbf{1}\|_{\infty} \leq \|Bw^* - \mathbf{1}\|_2 \leq \exp \left(-\Theta \left(\frac{L^2}{M} \right) \right). \quad (34)$$

In view of (27) and Lemma 2, we have

$$\|w^*\|_{\infty} \leq \|w^*\|_2 \leq \frac{\|\mathbf{1}\|_2}{\sigma_{\min}(B)} \leq \exp(O(L)).$$

Recall the connection between u_j and w_j in (7). For $1 \leq j \leq L < \beta \log k$, we have $u_j = w_j \frac{j!}{(\beta \log k)^j} \leq \frac{w_j}{\beta \log k}$. Therefore,

$$\|u^*\|_{\infty} \leq \frac{\|w^*\|_{\infty}}{\beta \log k} \leq \frac{\exp(O(L))}{\beta \log k}. \quad (35)$$

Applying (34) and (35) to Proposition 1, we obtain

$$\mathbb{E}(\hat{C} - C)^2 \leq k^2 \exp \left(-\frac{2n}{k} - \Theta \left(\frac{\alpha^2 n \log k}{\beta k} \right) \right) + k e^{-n/k} + k \frac{\exp(O(\alpha \log k))}{(\beta \log k)^2} + k^{-(\beta - \alpha \log \frac{\beta}{\alpha} - 3)}.$$

Then the desired (33) holds as long as β is sufficiently large and α is sufficiently small. \square

2.4 Lagrange interpolating polynomials and Stirling numbers

When we sample at least a constant fraction of the urn, i.e., $n = \Omega(k)$, we can afford to choose α and β in (8) so that $L = M$ and B is an invertible matrix. We choose the coefficient $w = B^{-1}\mathbf{1}$ which is equivalent to applying *Lagrange interpolating polynomial* and achieves exact zero bias. To control the variance, we can follow the approach in Section 2.3 by using the bound on minimum singular value of the matrix B , which implies that the coefficients is $\exp(O(L))$ and yields a coarse upper bound $O(k \frac{\log k}{1 \vee \log \frac{\Delta^2}{k}})$ on the sample complexity. As previously announced in Table 1, this bound can be improved to $O(k \log \frac{\log k}{1 \vee \log \frac{\Delta^2}{k}})$ by a more careful analysis of the Lagrange interpolating polynomial coefficients expressed in terms of the Stirling numbers, which we introduce next.

The Stirling numbers of the first kind are defined as the coefficients of the falling factorial $(x)_n$ where

$$(x)_n = x(x-1)\dots(x-n+1) = \sum_{j=1}^n s(n, j)x^j.$$

Compared to the coefficients w expressed by the Lagrange interpolating polynomial:

$$\sum_{j=1}^M w_j x^j - 1 = -\frac{(1-xM)(2-xM)\dots(M-xM)}{M!},$$

we obtain a formula for the coefficients w in terms of the Stirling numbers:

$$w_j = \frac{(-1)^{M+1} M^j}{M!} s(M+1, j+1), \quad 1 \leq j \leq M.$$

Consequently, the coefficients of our estimator u_j are given by

$$u_j = (-1)^{M+1} \frac{j!}{M!} \left(\frac{k}{n}\right)^j s(M+1, j+1). \quad (36)$$

The precise asymptotics the Stirling number is rather complicated. In particular, the asymptotic formula of $s(n, m)$ as $n \rightarrow \infty$ for fixed m is given by [Jor47] and the uniform asymptotics over all m is obtained in [MW58] and [Tem93]. The following lemma (proved in Appendix C) is a coarse non-asymptotic version, which suffices for the purpose of constant-factor approximations of the sample complexity.

Lemma 4.

$$|s(n+1, m+1)| = n! \left(\Theta \left(\frac{1}{m} \left(1 \vee \log \frac{n}{m} \right) \right) \right)^m \quad (37)$$

We construct \hat{C} as in Proposition 1 using the coefficients u_j in (36) to achieve zero bias. The variance upper bound by the coefficients u is a direct consequence of the upper bound of Stirling numbers in Lemma 4. Then we obtain the following MSE:

Theorem 2 (Interpolation). *Assume the Poisson sampling model. If $n > \eta k$ for some sufficiently large constant η , then*

$$\mathbb{E}(\hat{C} - C)^2 \leq k e^{-\Theta(\frac{n}{k})} + k^{-0.5-3.5\frac{k}{n} \log \frac{k}{en}} + \begin{cases} k \exp\left(\frac{k^2 \log k}{n^2} e^{-\Theta(\frac{n}{k})}\right), & n \lesssim k \log \log k, \\ k \left(\Theta\left(\frac{k}{n}\right) \log \frac{k^2 \log k}{n^2} \right)^{2n/k}, & k \log \log k \lesssim n \lesssim k \sqrt{\log k}, \\ 0, & n \gtrsim k \sqrt{\log k}. \end{cases}$$

Proof. In Proposition 1, fix $\beta = 3.5$ and $\alpha = \frac{\beta k}{n}$ so that $L = M$. Our goal is to show an upper bound of

$$\max_{\lambda = \frac{n}{k}[M]} \mathbb{E}_{N \sim \text{Poi}(\lambda)}[u_N^2] = \max_{\lambda = \frac{n}{k}[M]} \sum_{j=1}^M u_j^2 e^{-\lambda} \frac{\lambda^j}{j!}. \quad (38)$$

The coefficients u_j are obtained from (36) and (37):

$$|u_j| = \left(\Theta \left(\frac{k}{n} \right) \left(1 \vee \log \frac{M}{j} \right) \right)^j, \quad 1 \leq j \leq M. \quad (39)$$

Consider $n \gtrsim k\sqrt{\log k}$ and thus $\frac{n}{k} \gtrsim M$. The maximum of each summand in (38) as a function of $\lambda \in \mathbb{R}$ occurs at $\lambda = j$. Since $j \leq \frac{n}{k}$, the maximum over $\lambda = \frac{n}{k}[M]$ is attained at $\lambda = \frac{n}{k}$. Then,

$$\max_{\lambda = \frac{n}{k}[M]} \mathbb{E}_{N \sim \text{Poi}(\lambda)}[u_N^2] = \mathbb{E}_{N \sim \text{Poi}(\frac{n}{k})}[u_N^2]. \quad (40)$$

In view of (39) and $j \geq 1$, we have $|u_j| \leq (\Theta(k/n) \log M)^j$. Then,

$$\begin{aligned} \mathbb{E}_{N \sim \text{Poi}(\frac{n}{k})}[u_N^2] &\leq \mathbb{E}_{N \sim \text{Poi}(\frac{n}{k})} \left(\Theta \left(\frac{k \log M}{n} \right)^2 \right)^N \\ &= \exp \left(\frac{n}{k} \left(\Theta \left(\frac{k \log M}{n} \right)^2 - 1 \right) \right) = e^{-\Theta(n/k)}, \end{aligned}$$

as long as $n \gtrsim k \log \log k$ and thus $\frac{k \log M}{n} \lesssim 1$. Therefore,

$$\max_{\lambda = \frac{n}{k}[M]} \mathbb{E}_{N \sim \text{Poi}(\lambda)}[u_N^2] \leq e^{-\Theta(n/k)}, \quad n \gtrsim k\sqrt{\log k}. \quad (41)$$

For $k \log \log k \lesssim n \lesssim k\sqrt{\log k}$, we apply the following upper bound:

$$\begin{aligned} &\max_{\lambda = \frac{n}{k}[M]} \mathbb{E}_{N \sim \text{Poi}(\lambda)}[u_N^2] \\ &= \max_{\lambda = \frac{n}{k}[M]} \mathbb{E}_{N \sim \text{Poi}(\lambda)}[u_N^2 \mathbf{1}_{\{N \geq n/k\}}] + \max_{\lambda = \frac{n}{k}[M]} \mathbb{E}_{N \sim \text{Poi}(\lambda)}[u_N^2 \mathbf{1}_{\{N < n/k\}}] \\ &\leq \max_{\frac{n}{k} \leq j \leq M} |u_j|^2 + e^{-\Theta(n/k)}. \end{aligned} \quad (42)$$

where the upper bound of the second addend is analogous to (40) and (41). Since $\Theta(\frac{k}{n}) < 1$, the right-hand side of (39) is decreasing with j when $j \geq M/e$. It suffices to consider $j \leq M/e$, when the maximum as a function of $j \in \mathbb{R}$ occurs at $j = Me^{-\Theta(n/k)}$. Since $Me^{-\Theta(n/k)} \leq \frac{n}{k}$ when $n \gtrsim k \log \log k$, the maximum over $\frac{n}{k} \leq j \leq M$ is attained at $j = \frac{n}{k}$. Applying (39) with $j = \frac{n}{k}$ to (42) yields

$$\max_{\lambda = \frac{n}{k}[M]} \mathbb{E}_{N \sim \text{Poi}(\lambda)}[u_N^2] \leq \left(\Theta \left(\frac{k}{n} \right) \log \frac{k^2 \log k}{n^2} \right)^{2n/k} + e^{-\Theta(n/k)}. \quad (43)$$

Now consider $k \lesssim n \lesssim k \log \log k$. We apply the upper bound of expectation by the maximum:

$$\max_{\lambda = \frac{n}{k}[M]} \mathbb{E}_{N \sim \text{Poi}(\lambda)}[u_N^2] \leq \max_{j \in [M]} u_j^2.$$

Since $\Theta(\frac{k}{n}) < 1$, the right-hand side of (39) is decreasing with j when $j \geq M/e$, so it suffices to consider $j \leq M/e$. Denoting $x = \log \frac{M}{j}$ and $\tau = \Theta(\frac{k}{n})$, in view of (39), we have $|u_j| \leq \exp(Me^{-x} \log(\tau x))$, which attains maximum at x^* satisfying $\frac{e^{1/x^*}}{x^*} = \tau$. Then,

$$|u_j| \leq \exp(Me^{-x^*} \log(\tau x^*)) = \exp(Me^{-x^*}/x^*) < \exp(M\tau e^{-1/\tau}).$$

where the last inequality is because of $\tau > \frac{1}{x^*}$. Therefore,

$$\max_{\lambda = \frac{n}{k} [M]} \mathbb{E}_{N \sim \text{Poi}(\lambda)}[u_N^2] \leq \exp\left(\frac{k^2 \log k}{n^2} e^{-\Theta(\frac{n}{k})}\right), \quad k \lesssim n \lesssim k \log \log k. \quad (44)$$

Applying the upper bounds in (41), (43) and (44) to Proposition 1 concludes the proof. \square

Remark 5. It is impossible to bridge the gap near $\Delta = \sqrt{k}$ in Table 1 using the technology of interpolating polynomials that aims at zero bias, since its worst-case variance is at least $k^{1+\Omega(1)}$ when $n = O(k)$. To see this, note that from the variance term given by (12) is

$$\sum_{p_i} \mathbb{E}_{N \sim \text{Poi}(np_i)}[u_N^2] = \sum_{p_i} \sum_{j=1}^L u_j^2 e^{-np_i} \frac{(np_i)^j}{j!}. \quad (45)$$

Consider the distribution $\text{Uniform}[n/j_0]$ with $j_0 = Le^{-2n/k} = \Omega(\log k)$, which corresponds to a urn where each of the n/j_0 colors appears equal number of times. By the formula of coefficient u_j in (36) and the characterization from Lemma 4, the $j = j_0$ term in the summation of (45) is of order $\frac{n}{j_0} (\frac{k}{n} \log \frac{M}{j_0})^{2j_0} = \frac{n}{j_0} 2^{2j_0}$, which is already $k^{1+\Omega(1)}$.

3 Optimality of the sample complexity

In this section we develop lower bounds of the sample complexity which certifies the optimality of estimators constructed in Section 2. We first give a brief overview of the lower bound in [CCMN00, Theorem 1], which gives the optimal sample complexity under the multiplicative error criterion. The lower bound argument boils down to consider two hypothesis: in the null hypothesis, the urn consists of only one color; in the alternative, the urn contains $2\Delta + 1$ distinct colors, where $k - 2\Delta$ balls share the same color as in the null hypothesis, and all other balls have distinct colors. These two scenarios are distinguished if and only if a second color appears in the samples, which typically requires $\Omega(k/\Delta)$ samples. This lower bound is optimal for estimating within a multiplicative factor of $\sqrt{\Delta}$, which, however, is too loose for additive error Δ .

In contrast, instead of testing whether the urn is monochromatic, our first lower bound is based on given by testing whether the urn is maximally colorful, that is, containing k distinct colors. The alternative contains $k - 2\Delta$ colors, and the numbers of balls of two different colors differ by at most one. In other words, the null hypothesis is the uniform distribution on $[k]$ and the alternative is close to uniform distribution with smaller support size. The sample complexity, which is shown in Theorem 3, gives the lower bound in Table 1 for $\Delta \leq \sqrt{k}$.

Theorem 3. *If $1 \leq \Delta \leq \frac{k}{2}$, then*

$$n^*(k, \Delta) \geq \Omega\left(\frac{k - 2\Delta}{\sqrt{k}}\right). \quad (46)$$

If $1 \leq \Delta < k/4$, then

$$n^*(k, \Delta) \geq \Omega\left(k \operatorname{arccosh}\left(1 + \frac{k}{4\Delta^2}\right)\right) \asymp \begin{cases} k \log(1 + \frac{k}{\Delta^2}), & \Delta \leq \sqrt{k}, \\ \frac{k^{3/2}}{\Delta}, & \Delta \geq \sqrt{k}. \end{cases} \quad (47)$$

Proof. Consider the following two hypotheses: The null hypothesis H_0 is an urn consisting of k distinct colors; The alternative H_1 consists of $k - 2\Delta$ distinct colors, and each color appears either $b_1 \triangleq \lfloor \frac{k}{k-2\Delta} \rfloor$ or $b_2 \triangleq \lceil \frac{k}{k-2\Delta} \rceil$ times. In terms of distributions, H_0 is uniform distribution $Q = (\frac{1}{k}, \dots, \frac{1}{k})$; H_1 is the closest perturbations from uniform distributions: randomly pick disjoint indices $I, J \subseteq [k]$ with cardinality $|I| = c_1$ and $|J| = c_2$, where c_1 and c_2 satisfy

$$\begin{aligned} (\text{number of colors}) \quad c_1 + c_2 &= k - 2\Delta, \\ (\text{number of balls}) \quad c_1 b_1 + c_2 b_2 &= k. \end{aligned}$$

Conditional on $\theta \triangleq (I, J)$, the distribution $P_\theta = (p_{\theta,1}, \dots, p_{\theta,k})$ is given by

$$p_\theta = \begin{cases} b_1/k, & i \in I, \\ b_2/k, & i \in J. \end{cases}$$

Put the uniform prior on the alternative. Denote the marginal distributions of the n samples $X = (X_1, \dots, X_n)$ under H_0 and H_1 by Q_X and P_X , respectively. Since the distinct colors in H_0 and H_1 are separated by 2Δ , to show that the sample complexity $n^*(k, \Delta) \geq n$, it suffices to show that no test can distinguish H_0 and H_1 reliably using n samples. A further sufficient condition is a bounded χ^2 divergence [Tsy09] that

$$\chi^2(P_X \| Q_X) \triangleq \int \frac{P_X^2}{Q_X} - 1 \leq O(1).$$

The remainder of this proof is devoted to upper bounds of the χ^2 divergence.

Since $P_{X|\theta} = P_\theta^{\otimes n}$ and $Q_X = Q^{\otimes n}$, we have

$$\begin{aligned} \chi^2(P_X \| Q_X) + 1 &= \int \frac{P_X^2}{Q_X} = \int \frac{(\mathbb{E}_\theta P_{X|\theta})(\mathbb{E}_{\theta'} P_{X|\theta'})}{Q_X} \\ &= \mathbb{E}_{\theta, \theta'} \int \frac{P_{X|\theta} P_{X|\theta'}}{Q_X} = \mathbb{E}_{\theta, \theta'} \left(\int \frac{P_\theta P_{\theta'}}{Q} \right)^n, \end{aligned}$$

where θ' is an independent copy of θ . By the definition of P_θ and Q ,

$$\int \frac{P_\theta P_{\theta'}}{Q} = \frac{b_1^2}{k} |I \cap I'| + \frac{b_2^2}{k} |J \cap J'| + \frac{b_1 b_2}{k} (|I \cap J'| + |J \cap I'|) = 1 + \sum_{i=1}^4 A_i, \quad (48)$$

where $A_1 \triangleq \frac{b_1^2}{k} (|I \cap I'| - \frac{c_1^2}{k})$, $A_2 \triangleq \frac{b_2^2}{k} (|J \cap J'| - \frac{c_2^2}{k})$, $A_3 = \frac{b_1 b_2}{k} (|I \cap J'| - \frac{c_1 c_2}{k})$, and $A_4 = \frac{b_1 b_2}{k} (|J \cap I'| - \frac{c_1 c_2}{k})$ are centered random variables. Applying $1 + x \leq e^x$ and Cauchy-Schwartz inequality, we obtain

$$\chi^2(P_X \| Q_X) + 1 = \mathbb{E}[e^{n \sum_{i=1}^4 A_i}] \leq \prod_{i=1}^4 (\mathbb{E}[e^{4n A_i}])^{\frac{1}{4}}. \quad (49)$$

Consider $(\mathbb{E}[e^{4n A_1}])^{\frac{1}{4}}$. Note that $|I \cap I'| \sim \text{Hypergeometric}(k, c_1, c_1)$, which is the distribution of the sum of c_1 samples drawn without replacement from a population of size k which consists of c_1 ones and $k - c_1$ zeros. By the convex stochastic dominance of the binomial over the hypergeometric

distribution [Hoe63, Theorem 4], for $X \sim \text{Binomial}(c_1, \frac{c_1}{k})$, we have

$$\begin{aligned} (\mathbb{E}[e^{4nA_1}])^{\frac{1}{4}} &\leq \left(\mathbb{E} \left[\exp \left(\frac{4nb_1^2}{k} (X - c_1^2/k) \right) \right] \right)^{\frac{1}{4}} \\ &\leq \exp \left(\frac{c_1^2}{4k} \left(\exp \left(\frac{4nb_1^2}{k} \right) - 1 - \frac{4nb_1^2}{k} \right) \right) \\ &\leq \exp \left(\frac{c_1^2}{4k} \left(\exp \left(\frac{4nb_2^2}{k} \right) - 1 - \frac{4nb_2^2}{k} \right) \right), \end{aligned} \quad (50)$$

where the last inequality follows from the fact that $x \mapsto e^x - 1 - x$ is increasing when $x > 0$. Analogous upper bounds for other terms of (49) yields

$$\begin{aligned} &\chi^2(P_X \| Q_X) + 1 \\ &\leq \exp \left(\frac{c_1^2 + c_2^2 + 2c_1c_2}{4k} \left(\exp \left(\frac{4nb_2^2}{k} \right) - 1 - \frac{4nb_2^2}{k} \right) \right) \\ &= \exp \left(\frac{(k - 2\Delta)^2}{4k} \left(\exp \left(\frac{4n}{k} \left\lceil \frac{k}{k - 2\Delta} \right\rceil^2 \right) - 1 - \frac{4n}{k} \left\lceil \frac{k}{k - 2\Delta} \right\rceil^2 \right) \right). \end{aligned} \quad (51)$$

If $k - 2\Delta \geq \sqrt{k}$, the upper bound (51) implies that $n^*(k, \Delta) \geq \Omega(\frac{k-2\Delta}{\sqrt{k}})$; if $k - 2\Delta \leq \sqrt{k}$, the lower bound is trivial since $\frac{k-2\Delta}{\sqrt{k}} \leq 1$.

Now we proved the refined estimate (47) for $1 \leq \Delta < k/4$, in which case $|I| = c_1 = k - 4\Delta$, $|J| = c_2 = 2\Delta$ and $b_1 = 1, b_2 = 2$. c_1 is close to k , the hypergeometric distribution is no longer close to the binomial distribution, and the upper bound in (50) yields loose lower bound on sample complexity. In this case the set $K \triangleq (I \cup J)^c$ has small cardinality $|K| = 2\Delta$. The equality in (48) can be equivalent represented in terms of J, J' and K, K' by

$$\int \frac{P_\theta P_{\theta'}}{Q} = 1 + \frac{|J \cap J'| + |K \cap K'| - |J \cap K'| - |K \cap J'|}{k}.$$

By analogous upper bound to (49) – (51), $\chi^2(P_X \| Q_X) + 1 \leq \prod_{i=1}^4 (\mathbb{E}[e^{4nB_i}])^{\frac{1}{4}}$, where $B_1 \triangleq \frac{1}{k}(|J \cap J'| - \frac{(2\Delta)^2}{k})$, $B_2 \triangleq \frac{1}{k}(|K \cap K'| - \frac{(2\Delta)^2}{k})$, $B_3 \triangleq -\frac{1}{k}(|J \cap K'| - \frac{(2\Delta)^2}{k})$, and $B_4 \triangleq -\frac{1}{k}(|K \cap J'| - \frac{(2\Delta)^2}{k})$. For $X \sim \text{Binomial}(2\Delta, \frac{2\Delta}{k})$, we have

$$(\mathbb{E}[e^{4nB_i}])^{\frac{1}{4}} \leq \left(\mathbb{E} \left[\exp \left(t \left(X - \frac{(2\Delta)^2}{k} \right) \right) \right] \right)^{1/4} \leq \exp \left(\frac{(2\Delta)^2}{4k} (e^t - 1 - t) \right).$$

with $t = \frac{4n}{k}$ for $i = 1, 2$ and $t = -\frac{4n}{k}$ for $i = 3, 4$. Therefore,

$$\begin{aligned} \chi^2(P_X \| Q_X) + 1 &\leq \exp \left(\frac{\Delta^2}{k} (2e^{4n/k} + 2e^{-4n/k} - 4) \right) \\ &= \exp \left(\frac{4\Delta^2}{k} (\cosh(4n/k) - 1) \right). \end{aligned} \quad (52)$$

The upper bound (52) yields the sample complexity $n^*(k, \Delta) \geq \Omega(k \operatorname{arccosh}(1 + \frac{k}{4\Delta^2}))$. \square

Now we establish another lower bound for the sample complexity of the Distinct Elements problem for sampling without replacement. Since we can simulate sampling with replacement from

samples obtained without replacement (see (53) for details), it is also a valid lower bound for $n^*(k, \Delta)$ defined in Definition 1. On the other hand, as observed in [RRSS09, Lemma 3.3] (see also [Val12, Lemma 5.14]), any estimator \hat{C} for the Distinct Elements problem with sampling without replacement leads to an estimator for the Support Size problem with slightly worse performance: Suppose we have n i.i.d. samples drawn from a distribution P whose minimum non-zero probability is at least $1/l$. Let \hat{C}_{seen} denote the number of distinct elements in these samples. Equivalently, these samples can be viewed as being generated in two steps: first, we draw k i.i.d. samples from P , whose realizations form an instance of a k -ball urn with \hat{C}_{seen} distinct colors; next, we draw n samples from this urn without replacement ($n \leq k$), which clearly are distributed according to $P^{\otimes n}$. Suppose \hat{C}_{seen} is close to the actual support size of P . Then applying any algorithm for the Distinct Elements problem to these n i.i.d. samples constitutes a good support size estimator. Lemma 5 formalizes this intuition.

Lemma 5. *Suppose an estimator \hat{C} takes n samples from a k -ball urn ($n \leq k$) without replacement and provides an estimation error of less than Δ with probability at least $1 - \delta$. Applying \hat{C} with n i.i.d. samples from any distribution P with minimum non-zero mass $1/l$ and support size $S(P)$, we have*

$$|\hat{C} - S(P)| \leq 2\Delta$$

with probability at least $1 - \delta - \binom{l}{\Delta} \left(1 - \frac{\Delta}{l}\right)^k$.

Proof. Suppose that we take k i.i.d. samples from $P = (p_1, p_2, \dots)$, which form a k -ball urn consisting of C distinct colors. By the union bound,

$$\mathbb{P}[|C - S(P)| \geq \Delta] \leq \sum_{\substack{I: |I|=\Delta, \\ p_i \geq \frac{1}{l}, i \in I}} \left(1 - \sum_{i \in I} p_i\right)^k \leq \binom{l}{\Delta} \left(1 - \frac{\Delta}{l}\right)^k.$$

Next we take n samples without replacement from this urn and apply the given estimator \hat{C} . By assumption, conditioned on any realization of the k -ball urn, $|\hat{C} - C| \leq \Delta$ with probability at least $1 - \delta$. Then $|\hat{C} - S(P)| \leq 2\Delta$ with probability at least $1 - \delta - \binom{l}{\Delta} \left(1 - \frac{\Delta}{l}\right)^k$. Marginally, these n samples are identically distributed as n i.i.d. samples from P . \square

Combining with the sample complexity of the Support Size problem in (1), Lemma 5 leads to the following lower bound for the Distinct Elements problem:

Theorem 4. *Fix a sufficiently small constant c . For any $1 \leq \Delta \leq ck$,*

$$n^*(k, \Delta) \geq \Omega\left(\frac{k}{\log k} \log \frac{k}{\Delta}\right).$$

The same lower bound holds for sampling without replacement.

Proof. By the lower bound of the support size estimation problem obtained in [WY15, Theorem 2], if $n \leq \frac{\alpha l}{\log l} \log^2 \frac{l}{2\Delta}$ and $2\Delta \leq c_0 l$ for some fixed constants $c_0 < \frac{1}{2}$ and α , then for any \hat{C} , there exists a distribution P with minimum non-zero mass $1/l$ such that $|\hat{C} - S(P)| \leq 2\Delta$ with probability at most 0.8. Applying Lemma 5 yields that, using n samples without replacement, no estimator can provide an estimation error of Δ with probability 0.9 for an arbitrary k -ball urn, provided $\binom{l}{\Delta} \left(1 - \frac{\Delta}{l}\right)^k \leq 0.1$. Consequently, as long as $2\Delta \leq c_0 l$ and $\binom{l}{\Delta} \left(1 - \frac{\Delta}{l}\right)^k \leq 0.1$, we have

$$n^*(k, \Delta) \geq \frac{\alpha l}{\log l} \log^2 \frac{l}{2\Delta}.$$

The desired lower bound follows from choosing $l \asymp \frac{k}{\log(k/\Delta)}$. \square

A Connections between various sampling models

As mentioned in Section 1.2, four popular sampling models have been introduced in the statistics literature: multinomial model, hypergeometric model, Bernoulli model, and Poisson model. The connections between those models are explained in details in this section, as well as relations between the respective sample complexities.

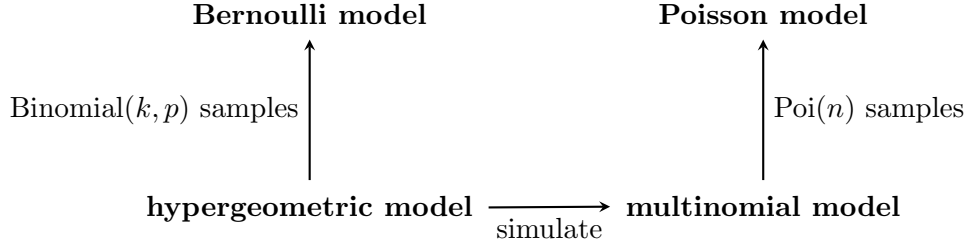


Figure 2: Relations between the four sampling models.

The connections between different models are illustrated in Fig. 2. Under the Poisson model, the sample size is a Poisson random variable; conditioned on the sample size, the samples are iid which is identical to the multinomial model. The same relation holds as the Bernoulli model to the hypergeometric model. Given samples (Y_1, \dots, Y_n) uniformly drawn from a k -ball urn without replacement (hypergeometric model), we can simulate (X_1, \dots, X_n) drawn with replacement (multinomial model) as follows: for each $i = 1, \dots, n$, let

$$X_i = \begin{cases} Y_i, & \text{with probability } 1 - \frac{i-1}{k}, \\ Y_m, & \text{with probability } \frac{1}{k}, \quad m \sim \text{Uniform}([i-1]). \end{cases} \quad (53)$$

In view of the connections in Fig. 2, any estimator constructed for one specific model can be adapted to another. The adaptation from multinomial to hypergeometric model is provided by the simulation in (53), and the other direction is given by Lemma 5 (without modifying the estimator). The following result provides a recipe for going between fixed and randomized sample size:

Lemma 6. *Let N be an \mathbb{N} -valued random variable.*

- (a) *Given any⁴ \hat{C} that uses n samples and succeeds with probability at least $1 - \delta$, there exists \hat{C}' using N samples that succeeds with probability at least $1 - \delta - \mathbb{P}[N < n]$.*
- (b) *Given any \tilde{C} using N samples that succeeds with probability at least $1 - \delta$, there exists \tilde{C}' that uses n samples and succeeds with probability at least $1 - \delta - \mathbb{P}[N > n]$.*

Proof. (a) Denote the samples by X_1, \dots, X_N . Following [RRSS09, Lemma 5.3(a)], define \hat{C}' as

$$\hat{C}' = \begin{cases} \hat{C}(X_1, \dots, X_n), & N \geq n, \\ 0, & N < n. \end{cases}$$

Then \hat{C}' succeeds as long as $N \geq n$ and \hat{C} succeeds, which has probability at least $1 - \delta - \mathbb{P}[N < n]$.

⁴More precisely, here and below \hat{C} is understood as a sequence of estimators indexed by the sample size $(X_1, \dots, X_n) \mapsto \hat{C}(X_1, \dots, X_n)$.

- (b) Denote the samples by X_1, \dots, X_n . Draw a random variable m from the distribution of N and define \tilde{C}' as

$$\tilde{C}' = \begin{cases} \tilde{C}(X_1, \dots, X_m), & m \leq n, \\ 0, & m > n. \end{cases}$$

The given estimator \tilde{C} fails with probability $\sum_{j \geq 0} \mathbb{P}[\tilde{C} \text{ fails} | N = j] \mathbb{P}[N = j] \leq \delta$. Consequently, $\sum_{j=0}^n \mathbb{P}[\tilde{C} \text{ fails} | N = j] \mathbb{P}[N = j] \leq \delta$. The estimator \tilde{C}' fails with probability at most

$$\sum_{j=0}^n \mathbb{P}[\tilde{C} \text{ fails} | m = j] \mathbb{P}[m = j] + \mathbb{P}[m > n] \leq \delta + \mathbb{P}[m > n].$$

We complete the proof. □

The adaptations of estimators between different sampling models imply the relations of the fundamental limits on the corresponding sample complexities. Extending Definition 1, let $n_M^*(k, \Delta, \delta)$, $n_H^*(k, \Delta, \delta)$, $n_B^*(k, \Delta, \delta)$, and $n_P^*(k, \Delta, \delta)$ be the minimum expected sample size under the multinomial, hypergeometric, Bernoulli, and Poisson sampling model, respectively, such that there exists an estimator \hat{C} satisfying $\mathbb{P}[|\hat{C} - C| \geq \Delta] \leq \delta$. Combining Chernoff bounds (see, e.g., [MU05, Theorem 4.4, 4.5 and 5.4]), we obtain Corollary 1, in which the connection between multinomial and Poisson models gives a rigorous justification of the assumption on the Poisson sampling model in Section 2.

Corollary 1. *The following relations hold:*

- n_H^* versus n_M^* :
 - (a) $n_H^*(k, \Delta, \delta) \leq n_M^*(k, \Delta, \delta)$;
 - (b) $n_H^*(k, \Delta, \delta) \leq n \Rightarrow n_M^*(k', 2\Delta, \delta + \binom{k'}{\Delta}(1 - \frac{\Delta}{k'})^k) \leq n$, for any $k' \in \mathbb{N}$. In particular, if δ is a constant, then we can choose $k' = \Theta(k / \log \frac{k}{\Delta})$.
- n_P^* versus n_M^* :
 - (c) $n_P^*(k, \Delta, \delta) \leq n \Rightarrow n_M^*(k, \Delta, \delta + (e/4)^n) \leq 2n$;
 - (d) $n_M^*(k, \Delta, \delta) \leq n \Rightarrow n_P^*(k, \Delta, \delta + (2/e)^n) \leq 2n$.
- n_B^* versus n_H^* :
 - (e) $n_B^*(k, \Delta, \delta) \leq n \Rightarrow n_H^*(k, \Delta, \delta + (e/4)^n) \leq 2n$;
 - (f) $n_H^*(k, \Delta, \delta) \leq n \Rightarrow n_B^*(k, \Delta, \delta + (2/e)^n) \leq 2n$.

B Correlation decay between fingerprints

Recall that the fingerprints are defined by $\Phi_j = \sum_i \mathbf{1}_{\{N_i=j\}}$, where N_i denotes the histogram of samples. Under the Poisson model, $N_i \stackrel{\text{ind}}{\sim} \text{Poi}(np_i)$. Then

$$\begin{aligned} \text{cov}(\Phi_j, \Phi_{j'}) &= - \sum_i \mathbb{P}[N_i = j] \mathbb{P}[N_i = j'], \quad j \neq j', \\ \text{var}[\Phi_j] &= \sum_i \mathbb{P}[N_i = j] (1 - \mathbb{P}[N_i = j]). \end{aligned}$$

The correlation coefficient between Φ_0 and Φ_j follows immediately:

$$\begin{aligned}
|\rho(\Phi_0, \Phi_j)| &= \sum_i \frac{\mathbb{P}[N_i = 0]\mathbb{P}[N_i = j]}{\sqrt{\sum_l \mathbb{P}[N_l = 0](1 - \mathbb{P}[N_l = 0]) \sum_l \mathbb{P}[N_l = j](1 - \mathbb{P}[N_l = j])}} \\
&\leq \sum_i \frac{\mathbb{P}[N_i = 0]\mathbb{P}[N_i = j]}{\sqrt{\mathbb{P}[N_i = 0](1 - \mathbb{P}[N_i = 0])\mathbb{P}[N_i = j](1 - \mathbb{P}[N_i = j])}} \\
&= \sum_i \sqrt{\frac{\mathbb{P}[N_i = 0]}{1 - \mathbb{P}[N_i = 0]} \frac{\mathbb{P}[N_i = j]}{1 - \mathbb{P}[N_i = j]}} = \sum_i \sqrt{\frac{e^{-\lambda_i}}{1 - e^{-\lambda_i}} \frac{\frac{e^{-\lambda_i} \lambda_i^j}{j!}}{1 - \frac{e^{-\lambda_i} \lambda_i^j}{j!}}}, \tag{54}
\end{aligned}$$

where $\lambda_i = np_i$. Note that $\max_{x>0} \frac{e^{-x}x^j}{j!} = \frac{e^{-j}j^j}{j!} \rightarrow 0$ as $j \rightarrow \infty$. Therefore, for any $x > 0$,

$$\frac{e^{-x}}{1 - e^{-x}} \frac{\frac{e^{-x}x^j}{j!}}{1 - \frac{e^{-x}x^j}{j!}} = \frac{1}{j!} \frac{e^{-2x}x^j}{1 - e^{-x}} (1 + o_j(1)), \tag{55}$$

where $o_j(1)$ is uniform as $j \rightarrow \infty$. Taking derivative, the function $x \mapsto \frac{e^{-2x}x^j}{1 - e^{-x}}$ on $x > 0$ is increasing if and only if $x + e^x(j - 2x) - j > 0$, and the maximum is attained at $x = j/2 + o_j(1)$. Therefore, applying $j! > (j/e)^j$,

$$\frac{1}{j!} \frac{e^{-2x}x^j}{1 - e^{-x}} \leq (1 + o_j(1))2^{-j}. \tag{56}$$

Combining (54) – (56), we conclude that

$$|\rho(\Phi_0, \Phi_j)| \leq k2^{-j/2}(1 + o_j(1)).$$

C Proof of auxiliary lemmas

Proof of Lemma 3. For any $z \in \mathbb{C}$, we can represent the forward difference in (20) as an integral:

$$\begin{aligned}
\Delta^m f(z) &= f(z + m) - \binom{m}{1} f(z + m - 1) + \cdots + (-1)^m f(z) \\
&= \int_{[0,1]^m} f^{(m)}(z + x_1 + \cdots + x_m) dx_1 \cdots dx_m.
\end{aligned}$$

Therefore,

$$|t_m(z)| = \left| \frac{1}{m!} \Delta^m p_m(z) \right| \leq \frac{1}{m!} \sup_{0 \leq \xi \leq m} |p_m^{(m)}(z + \xi)|. \tag{57}$$

Recall the definition of p_m in (21). Let $p_m(z) = \sum_{l=0}^{2m} a_l z^l$. Let $z(z-1)\cdots(z-m+1) = \sum_{i=0}^m b_i z^i$ and $(z-M)(z-M-1)\cdots(z-M-m+1) = \sum_{i=0}^m c_i z^i$. Expanding the product and collecting the coefficients yields a simple upper bound:

$$|b_i| \leq 2^m(m-1)^{m-i}, \quad |c_i| \leq 2^m(M+m-1)^{m-i} \leq 2^m(2M)^{m-i} \leq 2^{2m}M^{m-i}.$$

Since $\sum_{l=0}^{2m} a_l z^l = (\sum_{i=0}^m b_i z^i)(\sum_{j=0}^m c_j z^j)$, for $l \geq m$,

$$\begin{aligned}
|a_l| &= \left| \sum_{i=l-m}^m b_i c_{l-i} \right| \leq \sum_{i=l-m}^m 2^{3m}(m-1)^{m-i} M^{m-l+i} \\
&= 2^{3m} M^{2m-l} \sum_{i=l-m}^m \left(\frac{m-1}{M} \right)^{m-i} \leq m 2^{3m} M^{2m-l}.
\end{aligned}$$

Taking m -th derivative of p_m , we obtain

$$\begin{aligned}
|p_m^{(m)}(z)| &= \left| \sum_{j=0}^m a_{j+m} \frac{(j+m)!}{j!} z^j \right| \\
&\leq \sum_{j=0}^m |a_{j+m} M^j| \binom{m+j}{m} m! \left| \frac{z}{M} \right|^j \leq m 2^{3m} M^m m! (2e)^m \sum_{j=0}^m \left| \frac{z}{M} \right|^j \\
&\leq m^2 2^{6m} M^m m! \left(\frac{|z|}{M} \vee 1 \right)^m = m^2 2^{6m} m! (|z| \vee M)^m.
\end{aligned}$$

Then the desired (32) follows from (57). \square

Proof of Lemma 4. The following uniform asymptotic expansions of the Stirling numbers of the first kind have been obtain in [CRT00, Theorem 2]:

$$|s(n+1, m+1)| = \begin{cases} \frac{n!}{m!} (\log n + \gamma)^m (1 + o(1)), & 1 \leq m \leq \sqrt{\log n}, \\ \frac{\Gamma(n+1+R)}{\Gamma(R) R^{m+1} \sqrt{2\pi H}} (1 + o(1)), & \sqrt{\log n} \leq m \leq n - n^{1/3}, \\ \binom{n+1}{m+1} \left(\frac{m+1}{2} \right)^{n-m} (1 + o(1)), & n - n^{1/3} \leq m \leq n, \end{cases}$$

where γ is Euler's constant, R is the unique positive solution of $h'(x) = 0$ with $h(x) = \log \frac{\Gamma(x+n+1)}{\Gamma(x+1)x^m}$, $H = R^2 h''(R)$, and all $o(1)$ terms are uniform in m . In the following we consider each range separately and prove the non-asymptotic approximation in (37).

Case I. For $1 \leq m \leq \sqrt{\log n}$, Stirling's approximation gives

$$\frac{n!}{m!} (\log n + \gamma)^m = n! \left(\Theta \left(\frac{\log n}{m} \right) \right)^m.$$

Case II. For $n - n^{1/3} \leq m \leq n$,

$$\begin{aligned}
\binom{n+1}{m+1} \left(\frac{m+1}{2} \right)^{n-m} &= \frac{n!}{m!} \left(\Theta \left(\frac{m}{n-m} \right) \right)^{n-m} \\
&= n! \exp \left(m \left(\frac{n-m}{m} \log \left(\Theta \left(\frac{m}{n-m} \right) \right) - \log \Theta(m) \right) \right) \\
&= n! \left(\Theta \left(\frac{1}{m} \right) \right)^m.
\end{aligned}$$

Case III. For $\sqrt{\log n} \leq m \leq n - n^{1/3}$, note that $h(x) = \sum_{i=1}^n \log(x+i) - m \log x$, and thus $H = R^2 h''(R) = m - \sum_{i=1}^n \frac{R^2}{(R+i)^2} \leq m$. By [MW58, Lemma 4.1], $H = \omega(1)$ in this range. Hence,

$$|s(n+1, m+1)| = \frac{\Gamma(n+1+R)}{\Gamma(R) R^{m+1}} (\Theta(1))^m = \frac{n!}{R^m} \frac{\Gamma(n+1+R)}{n! \Gamma(R+1)} (\Theta(1))^m, \quad (58)$$

where R is the solution to $x(\frac{1}{x+1} + \dots + \frac{1}{x+n}) = m$. Bounding the sum by integrals, we have

$$R \log \left(1 + \frac{n}{R+1} \right) \leq m \leq R \log \left(1 + \frac{n}{R} \right).$$

If $\sqrt{\log n} \leq m \leq \frac{n}{e}$, then $R \asymp \frac{m}{\log(n/m)}$, and hence

$$1 \leq \frac{\Gamma(n+1+R)}{n!\Gamma(R+1)} \leq \left(O\left(\frac{n+R}{R} \right) \right)^R = \exp(O(m)).$$

In view of (58), we have $|s(n+1, m+1)| = \frac{n!}{(\Theta(R))^m}$, which is exactly (37) when $n \leq n/e$. If $n/e \leq m \leq n - n^{1/3}$, then $R \asymp \frac{n^2}{n-m}$, and

$$\begin{aligned} \frac{1}{R^m} \frac{\Gamma(n+1+R)}{n!\Gamma(R+1)} &= R^{-m} \left(\Theta\left(\frac{n+R}{n} \right) \right)^n \\ &= \exp\left(-m \log \Theta\left(\frac{n^2}{n-m} \right) + n \log \Theta\left(\frac{n}{n-m} \right) \right) \\ &= \exp\left(-m \log \Theta(n) + (n-m) \log \Theta\left(\frac{n}{n-m} \right) \right) \\ &= \exp(-m \log \Theta(n)). \end{aligned}$$

Combining (58) yields that $|s(n+1, m+1)| = n!(\Theta(\frac{1}{n}))^m$, which coincides with (37) since $n \asymp m$ is this range. \square

D Proof of results in Table 1

Below we explain how the sample complexity results summarized in Table 1 are obtained from various results in Section 2 and Section 3. The upper bounds are obtained from the worst-case MSE in Section 2 and the Markov inequality. In particular, the case of $\Delta \leq \sqrt{k}(\log k)^{-\delta}$ follows from the second and the third upper bounds of Theorem 2; the case of $\sqrt{k} \leq \Delta \leq k^{0.5+\delta}$ follows from the first upper bound of Theorem 2; the case of $k^{1-\delta} \leq \Delta \leq ck$ follows from Theorem 1. By monotonicity, we have the $O(k \log \log k)$ upper bound when $\sqrt{k}(\log k)^{-\delta} \leq \Delta \leq \sqrt{k}$, the $O(\frac{k}{\log k})$ upper bound when $\Delta \geq ck$, and the $O(k)$ upper bound when $k^{0.5+\delta} \leq \Delta \leq k^{1-\delta}$. The lower bound for $\Delta \leq \sqrt{k}$ follows from Theorem 3; the lower bound for $k^{0.5+\delta} \leq \Delta \leq ck$ follows from Theorem 4. These further implies the $\Omega(k)$ lower bound for $\sqrt{k} \leq \Delta \leq k^{0.5+\delta}$ by monotonicity.

References

- [Bec00] Bernhard Beckermann. The condition number of real Vandermonde, Krylov and positive definite Hankel matrices. *Numerische Mathematik*, 85(4):553–577, 2000.
- [BF93] John Bunge and M Fitzpatrick. Estimating the number of species: a review. *Journal of the American Statistical Association*, 88(421):364–373, 1993.
- [BYJK⁺02] Ziv Bar-Yossef, TS Jayram, Ravi Kumar, D Sivakumar, and Luca Trevisan. Counting distinct elements in a data stream. In *Proceedings of the 6th Randomization and Approximation Techniques in Computer Science*, pages 1–10. Springer-Verlag, 2002.
- [BYKS01] Ziv Bar-Yossef, Ravi Kumar, and D Sivakumar. Sampling algorithms: lower bounds and applications. In *Proceedings of the thirty-third annual ACM symposium on Theory of computing*, pages 266–275. ACM, 2001.

- [CCMN00] Moses Charikar, Surajit Chaudhuri, Rajeev Motwani, and Vivek Narasayya. Towards estimation error guarantees for distinct values. In *Proceedings of the nineteenth ACM SIGMOD-SIGACT-SIGART Symposium on Principles of Database Systems (PODS)*, pages 268–279. ACM, 2000.
- [CGR90] Antonio Córdova, Walter Gautschi, and Stephan Ruscheweyh. Vandermonde matrices on the circle: spectral properties and conditioning. *Numerische Mathematik*, 57(1):577–591, 1990.
- [CL92] Anne Chao and Shen-Ming Lee. Estimating the number of classes via sample coverage. *Journal of the American statistical Association*, 87(417):210–217, 1992.
- [CL99] Yang Chen and Nigel Lawrence. Small eigenvalues of large hankel matrices. *Journal of Physics A: Mathematical and General*, 32(42):7305, 1999.
- [CL11] T.T. Cai and M. G. Low. Testing composite hypotheses, Hermite polynomials and optimal estimation of a nonsmooth functional. *The Annals of Statistics*, 39(2):1012–1041, 2011.
- [CRT00] R Chelluri, LB Richmond, and NM Temme. Asymptotic estimates for generalized Stirling number. *Analysis-International Mathematical Journal of Analysis and its Application*, 20(1):1–14, 2000.
- [EPS01] Alfredo Eisinberg, Paolo Pugliese, and Nicola Salerno. Vandermonde matrices on integer nodes: the rectangular case. *Numerische Mathematik*, 87(4):663–674, 2001.
- [Est86] Warren W Esty. Estimation of the size of a coinage: A survey and comparison of methods. *The Numismatic Chronicle (1966-)*, pages 185–215, 1986.
- [ET76] B. Efron and R. Thisted. Estimating the number of unseen species: How many words did Shakespeare know? *Biometrika*, 63(3):435–447, 1976.
- [FCW43] Ronald Aylmer Fisher, A Steven Corbet, and Carrington B Williams. The relation between the number of species and the number of individuals in a random sample of an animal population. *The Journal of Animal Ecology*, pages 42–58, 1943.
- [Fer99] PJSG Ferreira. Super-resolution, the recovery of missing samples and Vandermonde matrices on the unit circle. In *Proceedings of the Workshop on Sampling Theory and Applications, Loen, Norway*, 1999.
- [FFGM07] Philippe Flajolet, Éric Fusy, Olivier Gandouet, and Frédéric Meunier. Hyperloglog: The analysis of a near-optimal cardinality estimation algorithm. In *In AofA07: Proceedings of the 2007 International Conference on Analysis of Algorithms*. Citeseer, 2007.
- [Fra78] Ove Frank. Estimation of the number of connected components in a graph by using a sampled subgraph. *Scandinavian Journal of Statistics*, pages 177–188, 1978.
- [Gau90] Walter Gautschi. How (un) stable are Vandermonde systems. *Asymptotic and computational analysis*, 124:193–210, 1990.
- [Goo49] Leo A Goodman. On the estimation of the number of classes in a population. *The Annals of Mathematical Statistics*, pages 572–579, 1949.

- [GS04] Alexander Goldenshluger and Vladimir Spokoiny. On the shape-from-moments problem and recovering edges from noisy Radon data. *Probability Theory and Related Fields*, 128(1):123–140, 2004.
- [Hil79] Bruce M Hill. Posterior moments of the number of species in a finite population and the posterior probability of finding a new species. *Journal of the American Statistical Association*, 74(367):668–673, 1979.
- [Hoe63] W. Hoeffding. Probability inequalities for sums of bounded random variables. *Journal of the American Statistical Association*, 58(301):13–30, Mar. 1963.
- [HOT88] Wen-Chi Hou, Gultekin Ozsoyoglu, and Baldeo K Taneja. Statistical estimators for relational algebra expressions. In *Proceedings of the seventh ACM SIGACT-SIGMOD-SIGART symposium on Principles of database systems*, pages 276–287. ACM, 1988.
- [Jor47] Charles Jordan. *Calculus of finite differences*. Chelsea, 1947.
- [KNW10] Daniel M Kane, Jelani Nelson, and David P Woodruff. An optimal algorithm for the distinct elements problem. In *Proceedings of the twenty-ninth ACM SIGMOD-SIGACT-SIGART symposium on Principles of database systems*, pages 41–52. ACM, 2010.
- [LNS99] Oleg Lepski, Arkady Nemirovski, and Vladimir Spokoiny. On estimation of the L_r norm of a regression function. *Probability theory and related fields*, 113(2):221–253, 1999.
- [Lo92] Shaw-Hwa Lo. From the species problem to a general coverage problem via a new interpretation. *The Annals of Statistics*, 20(2):1094–1109, 1992.
- [Moi15] Ankur Moitra. Super-resolution, extremal functions and the condition number of Vandermonde matrices. In *Proceedings of the Forty-Seventh Annual ACM on Symposium on Theory of Computing*, pages 821–830. ACM, 2015.
- [MU05] Michael Mitzenmacher and Eli Upfal. *Probability and computing: Randomized algorithms and probabilistic analysis*. Cambridge University Press, 2005.
- [MW58] L Moser and M Wyman. Asymptotic development of the Stirling numbers of the first kind. *Journal of the London Mathematical Society*, 1(2):133–146, 1958.
- [NS90] Jeffrey F Naughton and S Seshadri. On estimating the size of projections. In *International Conference on Database Theory*, pages 499–513. Springer, 1990.
- [NUS91] Arnold F Nikiforov, Vasilii B Uvarov, and Sergei K Suslov. *Classical orthogonal polynomials of a discrete variable*. Springer, 1991.
- [Pan03] Liam Paninski. Estimation of entropy and mutual information. *Neural Computation*, 15(6):1191–1253, 2003.
- [Pan04] Liam Paninski. Estimating entropy on m bins given fewer than m samples. *IEEE Transactions on Information Theory*, 50(9):2200–2203, 2004.
- [RRSS09] Sofya Raskhodnikova, Dana Ron, Amir Shpilka, and Adam Smith. Strong lower bounds for approximating distribution support size and the distinct elements problem. *SIAM Journal on Computing*, 39(3):813–842, 2009.

- [Sze75] G. Szegő. *Orthogonal polynomials*. American Mathematical Society, Providence, RI, 4th edition, 1975.
- [Tem93] Nico M Temme. Asymptotic estimates of Stirling numbers. *Studies in Applied Mathematics*, 89(3):233–243, 1993.
- [Tim63] Aleksandr Filippovich Timan. *Theory of approximation of functions of a real variable*. Pergamon Press, 1963.
- [Tod54] John Todd. The condition of the finite segments of the Hilbert matrix. *Contributions to the solution of systems of linear equations and the determination of eigenvalues*, 39:109–116, 1954.
- [Tsy09] A.B. Tsybakov. *Introduction to Nonparametric Estimation*. Springer Verlag, New York, NY, 2009.
- [Val11] Paul Valiant. Testing symmetric properties of distributions. *SIAM Journal on Computing*, 40(6):1927–1968, 2011.
- [Val12] Gregory Valiant. *Algorithmic Approaches to Statistical Questions*. PhD thesis, EECS Department, University of California, Berkeley, Sep 2012.
- [VV11a] Gregory Valiant and Paul Valiant. Estimating the unseen: an $n/\log(n)$ -sample estimator for entropy and support size, shown optimal via new CLTs. In *Proceedings of the 43rd annual ACM symposium on Theory of computing*, pages 685–694, 2011.
- [VV11b] Gregory Valiant and Paul Valiant. The power of linear estimators. In *Foundations of Computer Science (FOCS), 2011 IEEE 52nd Annual Symposium on*, pages 403–412. IEEE, 2011.
- [WY15] Yihong Wu and Pengkun Yang. Chebyshev polynomials, moment matching, and optimal estimation of the unseen. *arXiv:1504.01227*, 2015.
- [WY16] Yihong Wu and Pengkun Yang. Minimax rates of entropy estimation on large alphabets via best polynomial approximation. *IEEE Transactions on Information Theory*, 62(6):3702–3720, 2016.